# AN INVERSE PROBLEM FROM THE 2D-GROUNDWATER MODELLING

GOTTFRIED BRUCKNER, SYBILLE HANDROCK–MEYER, HARTMUT LANGMACH

Weierstrass Institute for Applied Analysis and Stochastics, Mohrenstrasse 39, D–10117 Berlin, Germany.

ABSTRACT. The paper is devoted to the inverse problem of identifying the coefficient in the main term of a quasilinear elliptic differential equation describing the filtration of groundwater. Experience suggests that the gradient of the piezometric head, i.e., Darcy's velocity, may have discontinuities and the transmissivity coefficient is a piecewise constant function.

For solving this problem we use a modification of a direct method of G. Vainikko. Starting with a weak formulation of the problem a suitable discretization is obtained by the method of minimal error. If necessary this method can be combined with Tikhonov regularization.

The main difficulty consists in generating distributed state observations from measurements of the ground–water level. For this step we propose an optimized data preparation procedure using additional information such as knowledge of the sought parameter values at some points and lower and upper bounds for the parameter.

Numerical tests show that locally sufficiently many measurements provide locally satisfactory results. Two numerical examples, one with simulated data and the other with real life data, are given.

## 1. INTRODUCTION

The two–dimensional steady flow in an isotropic and confined aquifer is governed, in general, by the quasilinear elliptic boundary value problem (cf. e.g. [5])

$$- \boldsymbol{\nabla} \cdot (a(\mathbf{x}, u)\boldsymbol{\nabla} u(\mathbf{x})) \; = \; f(\mathbf{x}) \quad \mathbf{x} \in \Omega \subset \mathbf{R}^2 \tag{1.1}$$

$$u(\mathbf{x}) \; = \; h(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega_1 \tag{1.2}$$

$$a(\mathbf{x}, u)\boldsymbol{\nabla} u(\mathbf{x}) \cdot \boldsymbol{\nu}(\mathbf{x}) \; = \; g(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega_2 = \partial\Omega \backslash \partial\Omega_1 \,, \tag{1.3}$$

where $\Omega$ is a bounded domain with piecewise smooth boundary and $\boldsymbol{\nu} = \boldsymbol{\nu}(\mathbf{x})$ is the outer unit normal on $\partial\Omega_2$. In the sequel, we confine ourselves to the special case that $\partial\Omega_1$ has positive Lebesgue measure and $h(\mathbf{x}) = h_0$. Physically, $u(\mathbf{x})$ can be interpreted as the groundwater level (piezometric head of ground water) in $\Omega$, and $a(\mathbf{x}, u)$ as transmissivity coefficient depending upon the space variable $\mathbf{x}$ and on the piezometric head $u(\mathbf{x})$. From this consideration it is clear that $a(\mathbf{x}, u) > 0$ for all admissible $\mathbf{x}$ and $u$. The function $f(\mathbf{x})$ characterizes sources or sinks in $\Omega$. The groundwater level on $\partial\Omega_1$ and the inflow or outflow through $\partial\Omega_2$ are denoted by $h_0$ and $g(\mathbf{x})$, respectively. The direct (forward) problem consists in the following:

Given $f, h_0, g, a$.   Find $u$.

For the well-posedness in the sense of Hadamard, (i.e. there exists a unique solution $u$ which continuously depends on the data $f, h_0, g, a$), of the direct problem (1.1)–(1.3) see [11]. Now let us formulate the inverse problem:

Given $f, h_0, g, u$.   Find $a$.

An inverse problem is ill–posed in general. Due to the lack of continuous dependence on the data (i.e. due to the lack of stability) difficulties arise when using noisy data.

Here we will be concerned with a stable reconstruction algorithm only, using Tikhonov regularization or self–regularization by discretization.

Let us briefly mention some relevant papers from the extensive literature concerning the inverse problem formulated above: Alessandrini [2] regularized the problem using singular perturbation theory which requires a high level of smoothness of the boundary and the data. In [14], [6], [15] the inverse problem is considered as a first order partial differential equation with respect to the unknown coefficient, where a high level of smoothness of the data has also to be assumed. A very fast procedure is obtained by the method of

Vainikko, where the inverse problem is transformed into a linear operator equation with a noncompact data dependent operator.

Hoffman and Sprekels [7], [8], [16] propose an adaptive method considering the steady state problem as an asymptotic limit of a suitable evolutionary process. The method needs the first derivatives of the data and has high numerical stability.

In this context the output least squares method is very often used, consisting of a minimum problem combined with Tikhonov regularization. Here no derivatives of the data are necessary, but the problem is nonlinear. Other methods using optimization procedures are the equation error method, the augmented Lagrange method [9] and the method of Lowe–Kohn [12]. In addition, let us refer to the papers [1], [4] as well as to the monograph [17].

To reduce the high computational expense of such methods, in this paper a direct inversion is proposed which is numerically cheap but very sensitive with respect to errors. Therefore, it is combined with an optimized data preparation procedure.

In our calculations we use what is basically a modification of Vainikko's method [18], [19], [20]. Starting with a weak formulation of the problem Vainikko's method consists of a finite element discretization of an operator equation in suitable Hilbert spaces, where the operator depends on the measured data. The considered projection method, the so–called method of least error, takes advantage of the simple form of the adjoint operator. The procedure is combined with Tikhonov regularization. This approach needs one measurement at each node.

In practice, however, only very few measurements are at our disposal so that data gained by interpolation are very erroneous and not in accordance with the a priori information on the coefficient. To counter these difficulties the method of Vainikko is combined with a method of "data smoothing" whose stabilizing effect consists of restricting the possible data set by a "smoothing" process. The goal of this method is an optimal utilization of the given information about both the coefficient and the data.

"New" data are sought, optimally fitting the "old" data and satisfying the discretized state equation with a certain tolerance, where the state equation is constructed using an a priori guess of the transmissivity. One gets a constrained minimization problem that is solved by the method of Lagrange multipliers and Newton's method. (Similar considerations in another context can be found in Parker's book [13].)

The paper can be understood as a continuation of [3] and is organized as follows. In Section 2 a short survey of the method is given. Section 3 deals with the data preparation. Finally, some numerical experiments are presented in Section 4.

## 2. THE METHOD

2.1. **Formulation of the problem and properties.** Let $\Omega \subset \mathbb{R}^d$ $(d \geq 2)$ be a bounded domain with piecewise smooth boundary $\partial\Omega$, where $\boldsymbol{\nu}(\mathbf{x})$ $(\mathbf{x} = (x_1, ..., x_d))$ is the outer unit normal on $\partial\Omega$. Furthermore, let $\partial\Omega_1 \subset \partial\Omega$ be a subset with a positive Lebesgue measure and $\partial\Omega_2 = \partial\Omega \backslash \partial\Omega_1$ be a relatively open subset having both, a piecewise smooth boundary on $\partial\Omega$. For a function $u \in W^{1,\infty}(\Omega)$ we define the real numbers

$$v_1 = \min_{\mathbf{x} \in \overline{\Omega}} u(\mathbf{x}), \qquad\qquad v_2 = \max_{\mathbf{x} \in \overline{\Omega}} u(\mathbf{x}),$$

and suppose that $v_1 < v_2$. In general the function $a(\mathbf{x}, u)$ is defined for all $\mathbf{x} \in \overline{\Omega}$ and all $u \in [v_1, v_2]$. The inverse problem describing the filtration of ground water in the domain $\Omega$ can be formulated in the following way:

Find the coefficient $a(\mathbf{x}, u) \in L^2(\Omega \times (v_1, v_2))$ such that

$$-\boldsymbol{\nabla} \cdot (a(\mathbf{x}, u)\boldsymbol{\nabla} u(\mathbf{x})) = f(\mathbf{x}) \quad \mathbf{x} \in \Omega \subset \mathbb{R}^d \tag{2.1}$$

$$u(\mathbf{x}) = 0 \qquad \mathbf{x} \in \partial\Omega_1 \tag{2.2}$$

$$a(\mathbf{x}, u)\boldsymbol{\nabla} u(\mathbf{x}) \cdot \boldsymbol{\nu}(\mathbf{x}) = g(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega_2 \,, \tag{2.3}$$

where $u \in W^{1,\infty}(\Omega)$, $f \in L^2(\Omega)$, $g \in L^2(\partial\Omega_2)$. Here $\partial\Omega_2 \subset \partial\Omega$ may be empty.

Physically, $u$ can be interpreted as the piezometric head of the groundwater in $\Omega$, the function $f$ characterizes the sources and sinks in $\Omega$ and the function $g$ describes the inflow and outflow through $\partial\Omega_2 \subset \partial\Omega$. The transmissivity coefficient $a$ is, physically, positive and piecewise continuous with possible discontinuities on some surfaces in $\Omega$.

Introducing the subspace

$$H^1(\Omega, \partial\Omega_1) = \{w \in H^1(\Omega) : w(\mathbf{x}) = 0 \quad \text{for } \mathbf{x} \in \partial\Omega_1\} \subset H^1(\Omega),$$

we can give the following weak formulation of the inverse problem (2.1)–(2.3): For given $u$ find $a \in L^2(\Omega \times (v_1, v_2))$ such that

$$\int_\Omega a\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w \, d\mathbf{x} = \int_\Omega fw \, d\mathbf{x} + \int_{\partial\Omega_2} gw \, dS \quad \text{for all } w \in H^1(\Omega, \partial\Omega_1) \,. \tag{2.4}$$

The problem (2.4) makes sense for $a \in L^2(\Omega \times (v_1, v_2))$ and $u \in W^{1,\infty}(\Omega)$.

We consider an auxiliary problem:

$$-\Delta\psi(\mathbf{x}) = f(\mathbf{x}) \quad \mathbf{x} \in \Omega \subset \mathbb{R}^d \tag{2.5}$$

$$\psi(\mathbf{x}) = 0 \qquad \mathbf{x} \in \partial\Omega_1 \tag{2.6}$$

$$\boldsymbol{\nabla}\psi(\mathbf{x}) \cdot \boldsymbol{\nu}(\mathbf{x}) = g(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega_2 \tag{2.7}$$

with the weak formulation

$$\int_\Omega \boldsymbol{\nabla}\psi \cdot \boldsymbol{\nabla} w \, d\mathbf{x} = \int_\Omega fw \, d\mathbf{x} + \int_{\partial\Omega_2} gw \, dS \quad \text{for all } w \in H^1(\Omega, \partial\Omega_1) \,. \tag{2.8}$$

From (2.4) and (2.8) we obtain

$$\int_\Omega a\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w \, d\mathbf{x} = \int_\Omega \boldsymbol{\nabla}\psi \cdot \boldsymbol{\nabla} w \, d\mathbf{x}. \tag{2.9}$$

Let $G$ be the space of gradients of functions $w \in H^1(\Omega, \partial\Omega_1)$:

$$G = G(\Omega, \partial\Omega_1) = \{\boldsymbol{\nabla} w : w \in H^1(\Omega, \partial\Omega_1)\} \subset (L^2(\Omega))^d \,.$$

Furthermore, using the orthoprojector

$$Q_G : (L^2(\Omega))^d \to G$$

we define an operator $T \in \mathcal{L}(L^2(\Omega \times (v_1, v_2)), G)$ by

$$Ta = Q_G(a\boldsymbol{\nabla} u), \quad a \in L^2(\Omega \times (v_1, v_2)) \tag{2.10}$$

and consider the operator equation

$$Ta = \nabla\psi, \tag{2.11}$$

where $\psi \in H^1(\Omega, \partial\Omega_1)$ is the solution of the direct problem (2.5)–(2.7). Using $Q_G\nabla\psi = \nabla\psi$ it is easily seen that the problem (2.4) is equivalent to the operator equation (2.11). If $\partial\Omega_2 \neq \partial\Omega$, then the direct problem (2.5)–(2.7) is uniquely solvable.

For general coefficients $a(\mathbf{x}, u)$ the equation (2.9) should be used for the discretization of the problem. But, sometimes the adjoint operator has a very simple form which can more conveniently be used for discretization.

Let us list some important special cases. From (2.11) we obtain

$$\langle Ta, \nabla w \rangle = \langle \nabla\psi, \nabla w \rangle = \langle a, T^*\nabla w \rangle \quad \text{for all} \quad \nabla w \in G, \tag{2.12}$$

where $\langle \cdot, \cdot \rangle$ denotes the scalar product in $L^2$. We restrict the function $a(\mathbf{x}, u)$ from $D(a) = \overline{\Omega} \times [v_1, v_2]$ to the graph $\{(\mathbf{x}, u(\mathbf{x})) \mid \mathbf{x} \in D(u) = \overline{\Omega}\}$ of the function $u$ and denote the restriction by $a(\mathbf{x}, u(\mathbf{x}))$.

1. We assume that $a(\mathbf{x}, u(\mathbf{x}))$ is a function depending only on the first variable $\mathbf{x}$, i.e. we identify $a(\mathbf{x}, u(\mathbf{x})) = a(\mathbf{x}) \in L^2(\Omega)$.
   Analogously to (2.10), we define $T_1 a = (Q_G a \nabla u)$ and using (2.12) we obtain

   $$T_1^* \nabla w = \nabla u \cdot \nabla w. \tag{2.13}$$

   This is the operator used in the method of Vainikko [18], [19], [20]. Vainikko was the first to recognize the advantages of applying the adjoint operator to the determination of coefficients in the form $a(\mathbf{x})$.

2. We assume that $a(\mathbf{x}, u(\mathbf{x}))$ is a function depending only on the second variable $u(\mathbf{x})$, i.e. we identify $a = a(\mathbf{x}, u(\mathbf{x})) = a(u(\mathbf{x})) \in L^2(\Omega)$.
   The transformation

   $$v(\mathbf{x}) = \int_{v_1}^{u(\mathbf{x})} a(s)\, ds \tag{2.14}$$

   yields

   $$\nabla v = a \nabla u, \tag{2.15}$$

   and

   $$-\Delta v(\mathbf{x}) = f(\mathbf{x}) \quad \mathbf{x} \in \Omega \tag{2.16}$$
   $$v(\mathbf{x}) = v_0 \quad \mathbf{x} \in \partial\Omega_1 \tag{2.17}$$
   $$\nabla v(\mathbf{x}) \cdot \boldsymbol{\nu}(\mathbf{x}) = g(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega_2 \tag{2.18}$$

   Now, the determination of the coefficient $a = a(u(\mathbf{x}))$ can be carried out in two steps:
   (a) Find $v$ from the direct problem (2.16)–(2.18).
   (b) Solve the integral equation (2.14)
      If $|\nabla u| \geq c > 0$, then $a(\mathbf{x}, u(\mathbf{x}))$ can be calculated via the formula

      $$a(u(\mathbf{x})) = \frac{\nabla v \cdot \nabla u}{|\nabla u|^2}$$

   which is obtained from (2.15).

3. We assume that $a(\mathbf{x}, u(\mathbf{x}))$ has the form
$$a(\mathbf{x}, u(\mathbf{x})) = a_1(\mathbf{x})\, a_2(u(\mathbf{x})), \qquad \mathbf{x} \in \overline{\Omega},$$
where $a_2(u(\mathbf{x}))$ is a known continuous function, i.e. we identify $a_1(\mathbf{x}) \in L^2(\Omega)$. Moreover, we assume that $a_1(\mathbf{x}) > 0$ and $a_2(u(\mathbf{x})) > 0$ for all $\mathbf{x} \in \overline{\Omega}$. Now, we define an operator $T_{a_2} \in \mathcal{L}(L^2(\Omega), G)$ by
$$T_{a_2}\, a_1 = Q_G(a_1\, a_2\, \boldsymbol{\nabla} u), \qquad a_1 \in L^2(\Omega). \tag{2.19}$$
Comparing (2.10) with (2.19) we see that $T_{a_2}\, a_1 = T\, a$ for $a = a_1\, a_2 \in L^2(\Omega)$. From (2.19) we conclude
$$\begin{aligned} \langle T_{a_2}\, a_1, \boldsymbol{\nabla} w \rangle &= \langle Q_G(a_1\, a_2\, \boldsymbol{\nabla} u), \boldsymbol{\nabla} w \rangle = \langle a_1\, a_2\, \boldsymbol{\nabla} u, Q_G\, \boldsymbol{\nabla} w \rangle \\ &= \langle a_1\, a_2\, \boldsymbol{\nabla} u, \boldsymbol{\nabla} w \rangle = \langle a_1,\, a_2\, \boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w \rangle \\ &= \langle a_1, T_{a_2}^*\, \boldsymbol{\nabla} w \rangle = \langle \boldsymbol{\nabla} \psi, \boldsymbol{\nabla} w \rangle \quad \text{for all} \quad \boldsymbol{\nabla} w \in G, \end{aligned}$$
where the adjoint operator $T_{a_2}^* \in \mathcal{L}L^2(\Omega), G)$ has the form
$$T_{a_2}^*\, \boldsymbol{\nabla} w = a_2(u)\, \boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w \quad \boldsymbol{\nabla} w \in G.$$
Using (2.13) we have with $K = K^* = a_2\, I : L^2(\Omega) \to L^2(\Omega)$
$$T_{a_2}^*\, \boldsymbol{\nabla} w = K\, T_1^*\, \boldsymbol{\nabla} w \quad \boldsymbol{\nabla} w \in G. \tag{2.20}$$
The operators $T_{a_2}^*$ and $T_{a_2}$ have the following properties:

1. Let be $d \geq 2$. Then the range $R(T_{a_2}^*) \subset L^2(\Omega)$ is nonclosed in $L^2(\Omega)$ even if $\mid \boldsymbol{\nabla} u \mid \geq c_1 > 0$ and $\mid a_2(u(\mathbf{x})) \mid \geq c_2 > 0$ in $\Omega$.
   Indeed, the nonclosedness of the range $R(T_1^*)$ of the operator $T_1^*$ is shown in [19]. Using (2.20) we see that $R(T_1^*) \subset D(K) = L^2(\Omega)$, $N(K) = \{0\}$ and $K$ is bounded in $L^2(\Omega)$. Then from $\overline{R(T_1^*)} \neq R(T_1^*)$ it follows that $\overline{R(KT_1^*)} \neq R(KT_1^*)$ ([10] $ 10).
2. $T_{a_2}^*$ is noncompact.
   $T_1^*$ is noncompact as a multiplication operator in the pair of spaces $(G, L^2(\Omega))$. Then the product $KT_1^*$, where $\mid a_2(u(\mathbf{x})) \mid \geq c_2 > 0$ in $\Omega$ is noncompact too.
3. The operator $T_{a_2}$ has a nonclosed range $R(T_{a_2}^*) \subset G$ and is also noncompact.
4. The problem (2.11) with the operator $T_1$ or $T_{a_2}$ is ill-posed.

## 2.2. Discretization and implementation.
As we mentioned above, for the discretization we use the simple form of the adjoint operator $T^*$, where for $T^*$ we take $T_1^*$ or $T_{a_2}^*$. The discretization is carried out by the method of minimal error, which is a special projection method.

Consider finite dimensional subspaces $S_h \subset H^1(\Omega, \partial\Omega_1)$ with the usual admissibility properties and take
$$G_h = \boldsymbol{\nabla} S_h \subset G, \quad T^*\boldsymbol{\nabla} S_h \subset L_2(\Omega)$$
as test and trial spaces, respectively. Then from (2.12)
$$\langle a_h, T^*\boldsymbol{\nabla} v_h \rangle = \langle \boldsymbol{\nabla} \psi, \boldsymbol{\nabla} v_h \rangle \quad \text{for all} \quad v_h \in S_h, \tag{2.21}$$
where $a_h = \boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} \hat{v}_h$, $\hat{v}_h \in S_h$, and
$$\|a_h - a\|_{L_2} = \min_{v_h \in S_h} \|\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} v_h - a\|_{L_2}.$$

Problem (2.21) has a unique solution $a_h$ and
$$\|a_h - a\|_{L_2} \to 0 \quad \text{as} \quad h \to 0.$$

Here $a$ is a minimal norm solution of (2.4).

The implementation is performed if $d = 2$, $\quad \mathbf{x} = (x_1, x_2)$, and for two types of coefficients:

1. $a(\mathbf{x}, u(\mathbf{x})) = a(\mathbf{x}) \quad \mathbf{x} \in \overline{\Omega}$,
2. $a(\mathbf{x}, u(\mathbf{x})) = a_1(\mathbf{x})(u(\mathbf{x}) - u_0(\mathbf{x})) \quad \mathbf{x} \in \overline{\Omega}$, where $a_1(\mathbf{x})$ denotes the transmissivity, $u(\mathbf{x})$ is the groundwater level and $u_0(\mathbf{x})$ is the lower bound of the aquifer. In particular, in our case, where the two-dimensional unconfined groundwater flow is considered, the coefficient is of this kind.

First we derive a linear equation system for the determination of the coefficient in the form $a(\mathbf{x})$ from the considerations in the special case 1 of Section 2.1. The linear equation system for the determination of the coefficient in the form $a_1(\mathbf{x})(u(\mathbf{x}) - u_0(\mathbf{x}))$ follows from the the special case 3 just there.

2.2.1. *Coefficients of the form $a = a(\mathbf{x})$.* Let $\Omega$ be a polygonal bounded domain and for a fixed discretization level $h$ let $\mathcal{T}_h$ be a regular triangulation, where

$$\overline{\Omega} = \bigcup_{E \in \mathcal{T}_h} \overline{E}.$$

Denote by $\mathcal{N} = \{\mathbf{P}_j\}_{j=1}^n$ the set of all nodes of the triangulation $\mathcal{T}_h$ that do not lie on the boundary $\partial\Omega_1$ and in the finite dimensional subspace $S_h \subset H^1(\Omega, \partial\Omega_1)$ choose a basis with linear base functions $\{w_j\}_{j=1}^n$ with $w_j = 0$ on $\partial\Omega_1$ and $w_j(\mathbf{P}_i) = \delta_{ij}$, $1 \leq i, j \leq n$.

Let us assume that the coefficient $a(\mathbf{x})$ is constant on each element (triangle) $E \in \mathcal{T}_h$ and the discretized coefficient $a_h$ can be represented as the vector

$$\mathbf{a} = (a^E)_{E \in \mathcal{T}_h}.$$

Then for the **direct problem**, where

$$u = \sum_{1 \leq j \leq n} u_j w_j$$

is to be determined, the linear system

$$\sum_{1 \leq j \leq n} L_{ij}[\mathbf{a}]u_j = d_i, \quad 1 \leq i \leq n, \tag{2.22}$$

where

$$L_{ij}[\mathbf{a}] \stackrel{def}{=} \sum_E a^E \int_E \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i \, d\mathbf{x}, \tag{2.23}$$

has to be solved ($a^E$ and $d_i$ are given). The values of $u$ on the boundary $\partial\Omega_1$ are already known as $u(\mathbf{x}) = 0$ on $\partial\Omega_1$.

For the **inverse problem** the linear system

$$\sum_{1 \leq j \leq n} M_{ij}[u]c_j = d_j, \quad 1 \leq i \leq n, \tag{2.24}$$

has to be solved, where

$$M_{ij}[u] \stackrel{def}{=} \sum_E \int_E (\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w_j)(\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w_i) \, d\mathbf{x}$$

(here $\boldsymbol{\nabla} u$ and $d_i$ are known). Then $a_h$ will be found from

$$a_h = \sum_{1 \le j \le n} c_j \boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w_j. \tag{2.25}$$

In combination with Tikhonov regularization (2.24) reads as

$$\sum_{1 \le j \le n} (\alpha L_{ij} + M_{ij}[u]) c_j^\alpha = d_j, \quad 1 \le i \le n, \tag{2.26}$$

where

$$L_{ij} \stackrel{def}{=} \sum_E \int_E \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i \, d\mathbf{x} = L_{ij}[\mathbf{1}].$$

It is clear that this method of Vainikko will work well when $\boldsymbol{\nabla} u$ has sufficiently good properties. If the matrix $(M_{ij}[u])_{i,j}$ in (2.24) is ill–conditioned, Tikhonov regularization (2.26) with a not too small $\alpha$ may produce results. However, if $\alpha$ is chosen too large the computed coefficient

$$a_h^\alpha = \sum_{1 \le j \le n} c_j^\alpha \boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w_j$$

cannot be interpreted as a solution to the inverse problem.

**Remark 2.1.** *The matrices* $\mathbf{L}[\mathbf{a}] = (L_{ij}[\mathbf{a}])_{i,j}$ *and* $\mathbf{M}[u] = (M_{ij}[u])_{i,j}$ *can be easily constructed using the coefficients*

$$L_{ij}^E = \int_E \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i \, d\mathbf{x} .$$

*Since* $\boldsymbol{\nabla} w_i \, (1 \le i \le n)$ *is constant on each element (i.e. triangle) $E$, we have*

$$L_{ij}^E = meas(E) \, \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i,$$

$$L_{ij}[\mathbf{a}] = \sum_E a^E \, L_{ij}^E,$$

$$M_{ji}[u] = M_{ij}[u] = \sum_E \Big( \sum_{\substack{k \in \mathcal{N} \\ 1 \le k \le n}} u_k \, L_{ki}^E \Big) \Big( \sum_{\substack{l \in \mathcal{N} \\ 1 \le l \le n}} u_l \, L_{lj}^E \Big) \frac{1}{meas(E)} .$$

$\square$

**Remark 2.2.** *If the triangle $E$ (of the triangulation $\mathcal{T}_h$ ) has no obtuse angles we have the well-known properties*

$$L_{ij}^E = L_{ji}^E \le 0, \quad i \ne j,$$

$$L_{ii}^E > 0,$$

$$\sum_{\substack{j \in \mathcal{N} \\ 1 \le j \le n}} L_{ij}^E = 0 \quad if \quad \overline{E} \cap \partial \Omega_1 = \emptyset.$$

In summary, choosing a basis $\{w_i\}_{i=1}^n$, setting in (2.21) $T^* = T_1^*$ and using (2.25) we obtain a linear equation system of the form

$$\mathbf{M}\,\mathbf{c} = \mathbf{d},$$

where $\mathbf{M} = (m_{ij})$ is a $n \times n$–matrix with the elements

$$m_{ij} = \langle T_1^* \boldsymbol{\nabla} w_j, T_1^* \boldsymbol{\nabla} w_i \rangle \quad i,j = 1,...,n,$$

and $\mathbf{c}$ as well as $\mathbf{d}$ are $n$–dimensional vectors with the components $c_1,...,c_n$

$$d_i = \langle \boldsymbol{\nabla}\psi, \boldsymbol{\nabla} w_i \rangle \quad i = 1,...,n, \tag{2.27}$$

respectively.

2.2.2. *Coefficients of the form* $a = a_1(\mathbf{x})(u(\mathbf{x}) - u_0(\mathbf{x}))$. It is easily seen that we can use the considerations of the special case 3 Subsection 2.1 if we set $a_2(u(\mathbf{x})) = u(\mathbf{x}) - u_0(\mathbf{x})$. Then (2.21) has the form

$$\langle a_{1h}, T_{a_2}^* \boldsymbol{\nabla} v_h \rangle = \langle \boldsymbol{\nabla}\psi, \boldsymbol{\nabla} v_h \rangle \quad \text{for all} \quad v_h \in S_h, \tag{2.28}$$

where

$$a_{1h} = \sum_{1 \le j \le n} c_{1j}\,\boldsymbol{\nabla} u \cdot \boldsymbol{\nabla} w_j.$$

Using as in Subsection 2.2.1 the basis $\{w_j\}_{j=1}^n$ we obtain from (2.28) a linear equation system in the form

$$\mathbf{N}\,\mathbf{c_1} = \mathbf{d},$$

where $\mathbf{N} = (n_{ij})$ is a $n \times n$-matrix with the elements

$$n_{ij} = \langle T_1^* \boldsymbol{\nabla} w_j, T_{a_2}^* \boldsymbol{\nabla} w_i \rangle \quad i,j = 1,...,n.$$

The n-dimensional vectors $\mathbf{c_1}$ and $\mathbf{d}$ have components $c_{11},...,c_{1n}$ and $d_i$ in the form (2.27).

**Remark 2.3.** *Let us compare the determination of a piecewise constant function* $a(\mathbf{x})$ *in the linear equation*

$$-\boldsymbol{\nabla} \cdot (a(\mathbf{x})\boldsymbol{\nabla} u(\mathbf{x})) = f(\mathbf{x}) \qquad \mathbf{x} \in \Omega \tag{2.29}$$

*with that of a piecewise constant function* $a_1(\mathbf{x})$ *in the quasilinear equation*

$$-\boldsymbol{\nabla} \cdot (a_1(\mathbf{x})(u(\mathbf{x}) - u_0(\mathbf{x}))\boldsymbol{\nabla} u(\mathbf{x})) = f(\mathbf{x}) \qquad \mathbf{x} \in \Omega. \tag{2.30}$$

*For the finite element discretization of (2.30) we use that* $a_1(\mathbf{x})$ *is piecewise constant, i.e.*

$$a_1(\mathbf{x}) = a_1^E \;\; for \;\; \mathbf{x} \in E,$$

*and suppose that* $u(\mathbf{x})$, $u_0(\mathbf{x})$ *are piecewise linear functions,*

$$
\begin{aligned}
u(\mathbf{x}) &= \sum_{1 \le k \le n} u_k w_k(\mathbf{x}), \\
u_0(\mathbf{x}) &= \sum_{1 \le k \le n} u_{0k} w_k(\mathbf{x}), \;\; \mathbf{x} \in \overline{\Omega}.
\end{aligned}
$$

*Recalling that $\boldsymbol{\nabla} w_k$ is constant on $E$, and $\int_E w_k(\mathbf{x})\, d\mathbf{x} = \frac{meas(E)}{d+1}$, if $P_k$ is a vertex of $E$, from (2.30) we obtain*

$$
\begin{aligned}
\int_\Omega f(\mathbf{x}) w_i(\mathbf{x})\, d\mathbf{x} &= \sum_E a_1^E \int_E \sum_{k, P_k \ vertex \ of \ E} (u_k - u_{0k}) w_k(\mathbf{x}) \sum_j u_j \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i\, d\mathbf{x} \\
&= \sum_E a_1^E \frac{meas(E)}{d+1} \sum_{k, P_k \ vertex \ of \ E} (u_k - u_{0k}) \sum_j u_j \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i.
\end{aligned}
$$

*On the other hand, from (2.29) we have*

$$
\int_\Omega f(\mathbf{x}) w_i(\mathbf{x})\, d\mathbf{x} = \sum_E a^E\, meas(E) \sum_j u_j \boldsymbol{\nabla} w_j \cdot \boldsymbol{\nabla} w_i
$$

*whence follows*

$$
a^E = a_1^E \frac{1}{d+1} \sum_{k, P_k \ vertex \ of \ E} (u_k - u_{0k}).
$$

*This means that $a_1^E$ can be determined from $a^E$ if $u(\mathbf{x}) > u_0(\mathbf{x})$ holds for every $\mathbf{x} \in \Omega$.*

$\square$

## 3. Description of the data smoothing procedure

3.1. **Preliminary remarks.** As in most inverse problems the influence of uncertain data is destructive to the inversion so that without regularization no useful result can be obtained.

In the problem considered here disturbances are caused on one hand by uncertain measurements of potential values and, on the other hand, by incomplete observations. To overcome the difficulties caused by noise Tikhonov regularization and the regularization by discretization had been proposed. Here, to avoid ambiguities caused by incomplete measurements, a so–called data smoothing procedure is considered. This procedure can be taken as some kind of regularization, where a well-behaved model is chosen which converges to the solution, if the noise (in this case the lack of measurements) tends to zero.

The inversion procedure of Vainikko, considered in this paper, needs one measurement at every node. The difficulty is that in practical tasks only very few measurements are at our disposal and, moreover these few measurements are not necessarily located at nodes in the domain.

The purpose of the data smoothing procedure is to construct a new data set suitable for the application of Vainikko's method. This suitable data set is to satisfy the state equation to a given tolerance and to have minimal distance from the measurement values. As well as the available measurements, a priori information is also of importance in the construction.

In what follows a matrix $\mathbf{B}$ relating the given information (measurements and a priori guesses) to the searched data set is defined and its properties are discussed. Then the minimum problem is formulated and solved.

The data smoothing procedure is described in detail. It is shown that the iterative application of this procedure decreases the distance between the calculated data and the measured ones.

Finally, some examples are considered and rules for the choice of procedure parameters are discussed.

The goal of these investigations is not to find the real permeability coefficient, which is impossible, because of the lack of measurements. Instead we attempt to use the given information in an optimal way, so as to be as helpful as possible in finding the true values.

**3.2. The matrix B.** First, let us recall some notation from 2.2. Let $\Omega$ be a bounded domain in $\mathbb{R}^2$ with a piecewise linear boundary, $\overline{\Omega} = \cup\overline{E}$ a fixed triangulation and $\mathcal{N}$ the set of nodes $\mathbf{P}_i$ $(i = 1, ..., n)$, being used in the Vainikko inversion. In addition, let us again consider the linear functions $w_i$ on $\Omega$ with the property $w_i(\mathbf{P}_j) = \delta_{ij}$.

To stress the correspondence between $w_i$ and $\mathbf{P}_i$, in what follows we shall write $w_{\mathbf{P}_i}$ instead of $w_i$. Then

$$w_{\mathbf{P}}(\mathbf{Q}) = \delta_{\mathbf{PQ}}, \quad \mathbf{P}, \mathbf{Q} \in \mathcal{N},$$

where

$$\delta_{\mathbf{PQ}} = \left\{ \begin{array}{ll} 1 & if \ \mathbf{P} = \mathbf{Q} \\ 0 & else. \end{array} \right.$$

Let

$$\mathcal{A} = \{\mathbf{a}, \ 0 < \alpha^E \leq a^E \leq \beta^E\},$$

be a set of admissible parameters, $\alpha^E$, $\beta^E$ given real numbers and $\mathbf{d}$ a fixed right hand side for the discretized direct problem (2.22). The vector $\mathbf{d}$, calculated from sources and boundary values, is assumed to be given exactly. Finally, let $\mathbf{a}_0$ be an a priori guess of the searched–for data set and

$$b_\lambda \quad \text{given measurements in} \quad \mathbf{y}_\lambda \in \Omega, \quad \lambda = 1, ..., m.$$

In what follows let us suppose

**Assumption 1:** Let $\mathbf{P} \in \mathcal{N}$ be fixed. If $w_{\mathbf{P}}(\mathbf{y}_\lambda) \geq w_{\mathbf{Q}}(\mathbf{y}_\lambda)$ for all $\mathbf{Q} \in \mathcal{N}$, then $\mathbf{y}_\lambda$ is uniquely determined.

Roughly speaking, assumption 1 means that the measurement points are more thinly distributed in $\Omega$ than the nodes.

Assumption 1 implies that a one–to–one relation can be defined between a subset $\mathcal{M} \subset \mathcal{N}$, $|\mathcal{M}| = m$ and the set of measurement points

$$\mathbf{P} \Longleftrightarrow \mathbf{x}_{\mathbf{P}}, \quad \mathbf{P} \in \mathcal{M}$$

with the property

$$w_{\mathbf{P}}(\mathbf{x}_{\mathbf{P}}) = \max_{\mathbf{Q} \in \mathcal{N}} w_{\mathbf{Q}}(\mathbf{x}_{\mathbf{P}}).$$

Indeed, let

$$\mathcal{M}_\lambda = \{\mathbf{P}, \ w_{\mathbf{P}}(\mathbf{y}_\lambda) \geq w_{\mathbf{Q}}(\mathbf{y}_\lambda) \quad \forall \mathbf{Q} \in \mathcal{N}\}.$$

The sets $\mathcal{M}_\lambda$ are non–empty and disjoint. Now let $\mathbf{P} \in \mathcal{M}_\lambda$ be arbitrarily chosen, and set $\mathbf{y}_\lambda = \mathbf{x}_{\mathbf{P}}$.

**Lemma 3.1.** *Let $V^n$ be the space of the functions*

$$u = \sum_{\mathbf{P} \in \mathcal{N}} u_{\mathbf{P}} w_{\mathbf{P}}, \quad u_{\mathbf{P}} \in I\!\!R,$$

*where $u(\mathbf{P}) = u_{\mathbf{P}}$. We consider the operator $\tilde{\mathbf{B}} : V^n \to V^n$*

$$\tilde{\mathbf{B}} w_{\mathbf{R}} = \begin{cases} w_{\mathbf{R}} + \sum_{\mathbf{P} \in \mathcal{M}} w_{\mathbf{R}}(\mathbf{x_P}) w_{\mathbf{P}}, & \mathbf{R} \in \mathcal{N} \backslash \mathcal{M} \\ \sum_{\mathbf{P} \in \mathcal{M}} w_{\mathbf{R}}(\mathbf{x_P}) w_{\mathbf{P}}, & \mathbf{R} \in \mathcal{M}, \end{cases}$$

*where the related matrix $\mathbf{B}$ mapping $I\!\!R^n$ to $I\!\!R^n$, has the entries*

$$B_{\mathbf{SR}} = \begin{cases} \delta_{\mathbf{RS}}, & \mathbf{S} \in \mathcal{N} \backslash \mathcal{M} \\ w_{\mathbf{R}}(\mathbf{x_S}), & \mathbf{S} \in \mathcal{M}. \end{cases} \tag{3.1}$$

*Then*

$$(\tilde{\mathbf{B}} u)_{\mathbf{P}} = \begin{cases} u(\mathbf{P}), & \mathbf{P} \in \mathcal{N} \backslash \mathcal{M} \\ u(\mathbf{x_P}), & \mathbf{P} \in \mathcal{M}. \end{cases}$$

*Proof.*

$$\tilde{\mathbf{B}} u = \sum_{\mathbf{R}} u(\mathbf{R}) \tilde{\mathbf{B}} w_{\mathbf{R}} = \sum_{\mathbf{R} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{R}) \tilde{\mathbf{B}} w_{\mathbf{R}} + \sum_{\mathbf{R} \in \mathcal{M}} u(\mathbf{R}) \tilde{\mathbf{B}} w_{\mathbf{R}}$$

$$= \sum_{\mathbf{R} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{R}) \left( w_{\mathbf{R}} + \sum_{\mathbf{P} \in \mathcal{M}} w_{\mathbf{R}}(\mathbf{x_P}) w_{\mathbf{P}} \right) + \sum_{\mathbf{R} \in \mathcal{M}} u(\mathbf{R}) \sum_{\mathbf{P} \in \mathcal{M}} w_{\mathbf{R}}(\mathbf{x_P}) w_{\mathbf{P}}$$

$$= \sum_{\mathbf{R} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{R}) w_{\mathbf{R}} + \sum_{\mathbf{P} \in \mathcal{M}} w_{\mathbf{P}} \left( \sum_{\mathbf{R} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{R}) w_{\mathbf{R}}(\mathbf{x_P}) + \sum_{\mathbf{R} \in \mathcal{M}} u(\mathbf{R}) w_{\mathbf{R}}(\mathbf{x_P}) \right)$$

$$= \sum_{\mathbf{P} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{P}) w_{\mathbf{P}} + \sum_{\mathbf{P} \in \mathcal{M}} \left( \sum_{\mathbf{R} \in \mathcal{N}} u(\mathbf{R}) w_{\mathbf{R}}(\mathbf{x_P}) \right) w_{\mathbf{P}}$$

$$= \sum_{\mathbf{P} \in \mathcal{N} \backslash \mathcal{M}} u(\mathbf{P}) w_{\mathbf{P}} + \sum_{\mathbf{P} \in \mathcal{M}} u(\mathbf{x_P}) w_{\mathbf{P}}.$$

$$\square$$

The matrix $\mathbf{B}$ is "nearly" a diagonal matrix, with $B_{\mathbf{RR}} \geq B_{\mathbf{RS}} \geq 0$, $B_{\mathbf{RR}} > 0$.

For the sequel, in addition to Assumption 1 let us suppose

**Assumption 2:** The matrix $\mathbf{B}$ is invertible.

**3.3. The minimum problem.** Starting from a given a priori guess $\mathbf{a}_0$, let $\mathbf{L}[\mathbf{a}_0]$ be constructed and let $\mathbf{u}[\mathbf{a}_0]$ be calculated by solving the direct problem

$$\mathbf{L}[\mathbf{a}_0] \mathbf{u}[\mathbf{a}_0] = \mathbf{d}.$$

Furthermore, we define the data vector $\bar{\mathbf{u}} = (\bar{u}_{\mathbf{P}}, \mathbf{P} \in \mathcal{N})$ as

$$\bar{u}_{\mathbf{P}} = \begin{cases} u[\mathbf{a}_0](\mathbf{P}); & \mathbf{P} \in \mathcal{N} \backslash \mathcal{M} \\ b_{\mathbf{P}}; & \mathbf{P} \in \mathcal{M}. \end{cases} \tag{3.2}$$

We have
$$\|\mathbf{B}u[\mathbf{a}_0] - \bar{\mathbf{u}}\|^2 = \sum_{\mathbf{P}\in\mathcal{M}} (u[\mathbf{a}_0](\mathbf{x_P}) - b_\mathbf{P})^2 = \|\mathbf{B}_m\mathbf{u}[\mathbf{a}_0] - \mathbf{b}\|_m^2 , \tag{3.3}$$

denoting by $\mathbf{B}_m$ the restriction of $\mathbf{B}$ to $\mathbb{R}^m$, the space of vectors $(u_\mathbf{P})_{\mathbf{P}\in\mathcal{M}}$.

Let us consider the minimum problem: Find $\hat{\mathbf{u}} \in V^n$ with the property

$$\|\mathbf{B}\hat{\mathbf{u}} - \bar{\mathbf{u}}\| = \min_{\{\mathbf{u}\in V^n,\, \|\mathbf{L}[\mathbf{a}_0]\mathbf{u}-\mathbf{d}\|\leq\delta\}} \|\mathbf{B}\mathbf{u} - \bar{\mathbf{u}}\| , \tag{3.4}$$

where $\delta > 0$ is given.

Let us assume that $\hat{\mathbf{u}}$ solves the state equation more exactly then $\mathbf{B}^{-1}\bar{\mathbf{u}}$. Then

$$0 \leq \delta \leq \|\mathbf{L}\mathbf{B}^{-1}\bar{\mathbf{u}} - \mathbf{d}\|.$$

**Theorem 3.1.** *Under the assumption 2 and if the operator $\mathbf{B}^T\mathbf{B} + \rho\mathbf{L}^T\mathbf{L}$ is invertible, the minimum problem (3.4) has exactly one solution*

$$\hat{\mathbf{u}} = (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d}),$$

*where $\mathbf{L} = \mathbf{L}[\mathbf{a}_0]$ and $\varrho$, $0 \leq \varrho \leq \infty$, is unique with the property*

$$\delta = \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}\mathbf{B}^T(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})\| .$$

*For $\varrho = \infty$ we have $\hat{\mathbf{u}} = \mathbf{u}[\mathbf{a}_0] = \mathbf{L}^{-1}\mathbf{d}$ and for $\varrho = 0$ we get $\hat{\mathbf{u}} = \mathbf{B}^{-1}\bar{\mathbf{u}}$.*

*Proof.* First, it is clear that (3.4) has a solution: Let $\{\mathbf{u}^k\}_k$ be a minimizing sequence, i.e.

$$\|\bar{\mathbf{u}} - \mathbf{B}\mathbf{u}^k\| \to \gamma := \inf_{\mathbf{u}\in V^n, \|\mathbf{L}\mathbf{u}-\mathbf{d}\|\leq\delta} \|\bar{\mathbf{u}} - \mathbf{B}\mathbf{u}\| .$$

Then $\|\mathbf{B}\mathbf{u}^k\| \leq C$, $\|\mathbf{u}^k\| \leq C$, for a subsequence $\{\mathbf{u}^{k_r}\}$ $\mathbf{u}^{k_r} \to \tilde{\mathbf{u}}$ as $r \to \infty$ and $\|\bar{\mathbf{u}} - \mathbf{B}\mathbf{u}^{k_r}\| \to \|\bar{\mathbf{u}} - \mathbf{B}\tilde{\mathbf{u}}\|$, $\|\mathbf{L}\mathbf{u}^{k_r} - \mathbf{d}\| \to \|\mathbf{L}\tilde{\mathbf{u}} - \mathbf{d}\|$ as $r \to \infty$, $\|\mathbf{L}\tilde{\mathbf{u}} - \mathbf{d}\| \leq \delta$, i.e. $\tilde{\mathbf{u}}$ is a solution of (3.4).

Additionally, if $\tilde{\mathbf{u}}$ solves (3.4) then $\tilde{\mathbf{u}}$ solves

$$\|\bar{\mathbf{u}} - \mathbf{B}\hat{\mathbf{u}}\| = \min_{\{\mathbf{u},\, \|\mathbf{L}\mathbf{u}-\mathbf{d}\|=\delta\}} \|\bar{\mathbf{u}} - \mathbf{B}\mathbf{u}\| . \tag{3.5}$$

Let $\tilde{\mathbf{u}}$ solve (3.4) with $\|\mathbf{L}\tilde{\mathbf{u}} - \mathbf{d}\| = \tilde{\delta} < \delta$. Let us show that $\|\mathbf{B}\tilde{\mathbf{u}} - \bar{\mathbf{u}}\| > 0$. Indeed, $\mathbf{B}\tilde{\mathbf{u}} = \bar{\mathbf{u}}$ means that $\tilde{\mathbf{u}} = \mathbf{B}^{-1}\bar{\mathbf{u}}$, and

$$\|\mathbf{L}\mathbf{B}^{-1}\bar{\mathbf{u}} - \mathbf{d}\| = \tilde{\delta} < \delta \leq \|\mathbf{L}\mathbf{B}^{-1}\bar{\mathbf{u}} - \mathbf{d}\|,$$

is a contradiction. Now consider

$$\mathbf{w} := \eta\tilde{\mathbf{u}} + (1 - \eta)\mathbf{B}^{-1}\bar{\mathbf{u}} \quad \text{for some} \quad \eta,\, 0 < \eta < 1 .$$

Then

$$\|\mathbf{B}\mathbf{w} - \bar{\mathbf{u}}\| = \eta\|\mathbf{B}\tilde{\mathbf{u}} - \bar{\mathbf{u}}\| < \|\mathbf{B}\tilde{\mathbf{u}} - \bar{\mathbf{u}}\|$$

for $\eta < 1$, and

$$\begin{aligned}\|\mathbf{L}\mathbf{w} - \mathbf{d}\| &= \|\mathbf{L}\tilde{\mathbf{u}} - \mathbf{d} + (1 - \eta)(\mathbf{L}\mathbf{B}^{-1}\bar{\mathbf{u}} - \mathbf{L}\tilde{\mathbf{u}})\| \\ &\leq \|\mathbf{L}\tilde{\mathbf{u}} - \mathbf{d}\| + (1 - \eta)C \leq \delta,\end{aligned}$$

which contradicts to the assumption that $\tilde{\mathbf{u}}$ solves (3.4). The proof of (3.5) is complete.

Now, let us solve (3.5). Consider the Lagrange function

$$\mathfrak{L}(\mathbf{u}, \alpha) = \|\mathbf{B}\mathbf{u} - \bar{\mathbf{u}}\|^2 + \alpha(\|\mathbf{L}\mathbf{u} - \mathbf{d}\|^2 - \delta^2) .$$

Necessary conditions for a minimum are

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = 0 \,, \quad \frac{\partial \mathcal{L}}{\partial \alpha} = 0 \,.$$

We have

$$\left(\frac{\partial \mathcal{L}}{\partial \mathbf{u}}, \underline{\boldsymbol{\delta}}\mathbf{u}\right) = \lim_{\lambda \to 0} \frac{1}{\lambda}(\|\mathbf{B}(\mathbf{u} + \lambda\underline{\boldsymbol{\delta}}\mathbf{u}) - \bar{\mathbf{u}}\|^2 - \|\mathbf{B}\mathbf{u} - \bar{\mathbf{u}}\|^2 + \alpha(\|\mathbf{L}(\mathbf{u} + \lambda\underline{\boldsymbol{\delta}}\mathbf{u}) - \mathbf{d}\|^2$$

$$-\|\mathbf{L}\mathbf{u} - \mathbf{d}\|^2)) = 2(\mathbf{B}^T(\mathbf{B}\mathbf{u} - \bar{\mathbf{u}}) + \alpha\mathbf{L}^T(\mathbf{L}\mathbf{u} - \mathbf{d}), \underline{\boldsymbol{\delta}}\mathbf{u}).$$

Then, we obtain from the necessary conditions

$$\mathbf{B}^T(\mathbf{B}\hat{\mathbf{u}} - \bar{\mathbf{u}}) + \varrho\mathbf{L}^T(\mathbf{L}\hat{\mathbf{u}} - \mathbf{d}) = 0 \,, \tag{3.6}$$

$$\delta = \|\mathbf{L}\hat{\mathbf{u}} - \mathbf{d}\| \,. \tag{3.7}$$

From the equations (3.6), (3.7) the pair $(\hat{\mathbf{u}}, \varrho)$ is uniquely determined, if $\mathbf{d}$, $\bar{\mathbf{u}}$, $\delta$ are given. Indeed, (3.6) implies

$$\hat{\mathbf{u}} = (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d}).$$

Then we have

$$\begin{aligned}
\delta &= \|\mathbf{L}\hat{\mathbf{u}} - \mathbf{d}\| \\
&= \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d}) - \mathbf{d}\| \\
&= \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d} - (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})\mathbf{L}^{-1}\mathbf{d})\| \\
&= \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d} - \mathbf{B}^T\mathbf{B}\mathbf{L}^{-1}\mathbf{d} - \varrho\mathbf{L}^T\mathbf{L}\mathbf{L}^{-1}\mathbf{d})\| \\
&= \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} - \mathbf{B}^T\mathbf{B}\mathbf{L}^{-1}\mathbf{d})\|
\end{aligned}$$

whence follows

$$\delta = \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}\mathbf{B}^T(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})\|. \tag{3.8}$$

The function

$$\theta(s) = \|\mathbf{L}(\mathbf{B}^T\mathbf{B} + s\mathbf{L}^T\mathbf{L})^{-1}\mathbf{B}^T(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})\|^2$$

is strictly decreasing for $0 \le s < \infty$. Indeed, introducing $\mathbf{z}(s)$ as the solution of

$$(\mathbf{B}^T\mathbf{B} + s\mathbf{L}^T\mathbf{L})\mathbf{z}(s) = \mathbf{B}^T(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})$$

we obtain

$$\theta'(s) = \frac{d}{ds}\langle \mathbf{L}\mathbf{z}(s), \mathbf{L}\mathbf{z}(s)\rangle = 2\langle \mathbf{L}\mathbf{z}'(s), \mathbf{L}\mathbf{z}(s)\rangle = 2\langle \mathbf{z}'(s), \mathbf{L}^T\mathbf{L}\mathbf{z}(s)\rangle,$$

where $\mathbf{z}'(s)$ is given as the solution of

$$(\mathbf{B}^T\mathbf{B} + s\mathbf{L}^T\mathbf{L})\mathbf{z}'(s) + \mathbf{L}^T\mathbf{L}\mathbf{z}(s) = 0 \,.$$

Then it follows

$$\theta'(s) = -2\langle(\mathbf{B}^T\mathbf{B} + s\mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^T\mathbf{L}\mathbf{z}(s), \mathbf{L}^T\mathbf{L}\mathbf{z}(s)\rangle \le 0$$

since $(\mathbf{B}^T\mathbf{B} + s\mathbf{L}^T\mathbf{L})^{-1}$ is positive definite. In addition it holds

$$\theta(0) = \|\mathbf{L}(\mathbf{B}^T\mathbf{B})^{-1}\mathbf{B}^T(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})\|^2 = \|\mathbf{L}\mathbf{B}^{-1}(\bar{\mathbf{u}} - \mathbf{B}\mathbf{L}^{-1}\mathbf{d})\|^2 = \|\mathbf{L}\mathbf{B}^{-1}\bar{\mathbf{u}} - \mathbf{d}\|^2 > 0.$$

Now, we obtain the uniqueness of $\varrho$ from (3.8). $\qquad\square$

**3.4. The procedure.** The purpose is the construction of a data set $\hat{\mathbf{u}}$ suitable for the Vainikko inversion, from the a priori guess $\mathbf{a}_0$ and the measurements $\mathbf{b}$, and to repeat this process if reasonable.

*Preparation:* Start with constructing $\mathbf{L}[\mathbf{a}_0] = \mathbf{L}$ from the given a priori guess $\mathbf{a}_0$, then calculate $\mathbf{u}[\mathbf{a}_0]$ by solving a direct problem. After defining the correspondence $\mathbf{x_P} \Longleftrightarrow \mathbf{P}$, $\mathbf{P} \in \mathcal{M}$ ( by calculating $w_{\mathbf{Q}}(\mathbf{y}_\lambda) \, \mathbf{Q} \in \mathcal{N}$, $\lambda = 1, ..., m$) construct the matrices $\mathbf{B}$ in (3.1) and $\mathbf{B}^T\mathbf{B}$, and build the vector $\bar{\mathbf{u}}$ from (3.2).

*Iteration:*

(I)       Choose $\varrho$,  $0 \leq \varrho \leq \infty$.

(II)      Calculate $\hat{\mathbf{u}} = (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}} + \varrho\mathbf{L}^T\mathbf{d})$.

(III)     Invert $\hat{\mathbf{u}}$ by the Vainikko method, the result is $\tilde{\mathbf{a}}$.

(IV)     Project $\tilde{\mathbf{a}}$ to the convex set $\mathcal{A}$, i.e. $\mathbf{a} = \mathbf{P}_\mathcal{A}\tilde{\mathbf{a}}$.

(V)      Apply a suitable stopping rule.

(VI)     Construct $\mathbf{L} = \mathbf{L}[\mathbf{a}]$.

(VII)    Construct $\bar{\mathbf{u}}$ :  $\bar{u}_\mathbf{P} = \begin{cases} \hat{u}_\mathbf{P}, & \mathbf{P} \in \mathcal{N}\backslash\mathcal{M} \\ b_\mathbf{P}, & \mathbf{P} \in \mathcal{M}, \end{cases}$

(VIII)  Go to (I).

**Theorem 3.2.** *Let the assumptions of Theorem 3.1 be satisfied and let $\hat{\mathbf{u}}_0 = \mathbf{u}[\mathbf{a}_0]$, $\hat{\mathbf{u}}_i = \hat{\mathbf{u}}$ in the i–th iteration of (II), $\bar{\mathbf{u}}_i = \bar{\mathbf{u}}$ in the i–th iteration of (VII). If $\varrho < \infty$ then the sequence $\hat{\mathbf{u}}_i$ has the property*

$$\|\mathbf{B}_m\hat{\mathbf{u}}_{i+1} - \mathbf{b}\|_m < \|\mathbf{B}_m\hat{\mathbf{u}}_i - \mathbf{b}\|_m, \quad i = 0, 1, 2, ..., \tag{3.9}$$

*(i.e. by iteration the fitting of the measurements will be improved).*

*Proof.* First let us prove the inequality

$$\|\mathbf{B}\hat{\mathbf{u}}_i - \bar{\mathbf{u}}_{i-1}\| < \|\mathbf{B}\hat{\mathbf{u}}_{i-1} - \bar{\mathbf{u}}_{i-1}\| \quad i = 1, 2, .... \tag{3.10}$$

Let $i$ be one of the numbers $1, 2, ...$ and

$$\mathbf{g} := \mathbf{B}\hat{\mathbf{u}}_{i-1} - \bar{\mathbf{u}}_{i-1} \neq 0. \tag{3.11}$$

For $\mathbf{g} = 0$, then we have a total fit of the measurements by $\hat{\mathbf{u}}_{i-1}$. In this case the inversion $\mathbf{a}_{i-1}$ of $\hat{\mathbf{u}}_{i-1}$ cannot be improved and has to be taken as solution.

Setting in (2.23) $\mathbf{L} = \mathbf{L}[\mathbf{a}_{i-1}]$ we have

$$\hat{\mathbf{u}}_i = (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}}_{i-1} + \varrho\mathbf{L}^T\mathbf{d}),$$

implying

$$\mathbf{B}\hat{\mathbf{u}}_i - \bar{\mathbf{u}}_{i-1} = \mathbf{B}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}(\mathbf{B}^T\bar{\mathbf{u}}_{i-1} + \varrho\mathbf{L}^T\mathbf{d} - (\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})\mathbf{B}^{-1}\bar{\mathbf{u}}_{i-1}).$$

Straightforward calculations yield

$$\mathbf{B}\hat{\mathbf{u}}_i - \bar{\mathbf{u}}_{i-1} = \mathbf{B}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}\varrho\mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}(\mathbf{B}\mathbf{L}^{-1}\mathbf{d} - \bar{\mathbf{u}}_{i-1}).$$

Since $\mathbf{a}_{i-1}$ corresponds to the inversion of $\hat{\mathbf{u}}_{i-1}$, under the assumption that the difference between $\mathbf{a}_{i-1}$ and $\tilde{\mathbf{a}}_{i-1}$ is small we have approximately

$$\mathbf{L}^{-1}\mathbf{d} = \hat{\mathbf{u}}_{i-1}.$$

Then we obtain

$$\mathbf{B}\hat{\mathbf{u}}_i - \bar{\mathbf{u}}_{i-1} = \mathbf{B}(\mathbf{B}^T\mathbf{B} + \varrho\mathbf{L}^T\mathbf{L})^{-1}\varrho\mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}\mathbf{g}$$
$$= \mathbf{B}(\varrho^{-1}\mathbf{B}^T\mathbf{B} + \mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}\mathbf{g}.$$

Now let us consider the continuous function

$$\psi(s) = \|\mathbf{B}\mathbf{z}(s)\|^2,$$

where $\mathbf{z}(s)$ is defined by

$$(s\mathbf{B}^T\mathbf{B} + \mathbf{L}^T\mathbf{L})\mathbf{z}(s) = \mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}\mathbf{g}. \tag{3.12}$$

Then

$$\psi'(s) = 2\langle\mathbf{B}\mathbf{z}'(s), \mathbf{B}\mathbf{z}(s)\rangle = 2\langle\mathbf{z}'(s), \mathbf{B}^T\mathbf{B}\mathbf{z}(s)\rangle,$$

where $\mathbf{z}'(s)$ can be determined by differentiating (3.12), i.e.

$$(s\mathbf{B}^T\mathbf{B} + \mathbf{L}^T\mathbf{L})\mathbf{z}'(s) + \mathbf{B}^T\mathbf{B}\mathbf{z}(s) = \mathbf{0}.$$

From this we have

$$\psi'(s) = -2\langle(s\mathbf{B}^T\mathbf{B} + \mathbf{L}^T\mathbf{L})^{-1}\mathbf{B}^T\mathbf{B}\mathbf{z}(s), \mathbf{B}^T\mathbf{B}\mathbf{z}(s)\rangle < 0$$

as $\mathbf{z}(s) \neq \mathbf{0}$ (c.f. (3.11)).

Furthermore,

$$\psi(0) = \|\mathbf{B}(\mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}\mathbf{g}\|^2 = \|\mathbf{g}\|^2 \neq 0,$$

$$\lim_{s\to\infty}\psi(s) = \lim_{s\to\infty} s^{-2}\|\mathbf{B}(\mathbf{B}^T\mathbf{B} + s^{-1}\mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^T\mathbf{L}\mathbf{B}^{-1}\mathbf{g}\|^2 = 0.$$

That means $\psi(s) < \psi(0)$ if $s > 0$, i.e. (3.10).

Now, let us prove (3.9). First of all we have (cf. (3.3))

$$\|\mathbf{B}\hat{\mathbf{u}}_j - \bar{\mathbf{u}}_j\|^2 = \sum_{\mathbf{P}\in\mathcal{M}}(\hat{u}_j(\mathbf{x_P}) - b_{\mathbf{P}})^2 = \|\mathbf{B}_m\hat{\mathbf{u}}_j - \mathbf{b}\|_m^2, \quad j = 0, 1, 2, .... \tag{3.13}$$

Since obviously

$$\|\mathbf{B}_m\hat{\mathbf{u}}_i - \mathbf{b}\|_m^2 \leq \|\mathbf{B}\hat{\mathbf{u}}_i - \bar{\mathbf{u}}_{i-1}\|^2$$

the assertion (3.9) follows from (3.10) and (3.13) for $j = i - 1$. □

Remarks on the choice of $\varrho$ and on the stopping rule (V) can be found below.

## 4. NUMERICAL EXAMPLES

4.1. **A numerical example with simulated data.** In the following numerical experiments the effect of the data preparation will be demonstrated. Let us consider a square domain $\Omega = \{\mathbf{x} \in \mathbf{R}^2 :| \mathbf{x} |< 1.1\}$ with impermeable upper and lower boundary, homogeneous Dirichlet conditions at the left boundary and inhomogeneous Neumann conditions at the right one. No sources and sinks are considered. The domain $\Omega$ is triangulated by $30 \times 30$ equidistant nodes.

*Data generation:* Suppose that we are given (cf. Fig. 1)

1. (from a geological a priori information) an a priori guess $\mathbf{a}_0$, i.e., more precisely, a lens $C \supset A$ of diminished (constant) transmissivity $a_{02} = 10^{-3}$ surrounded by an area $\Omega \backslash C$ of (constant) transmissivity $a_{01} = 10^{-5}$;

2. measurements in an area $B' \supset B$ at about 75% of all nodes in $B'$.

These measurements are simulated by the potential values resulting from an assumed reality, i.e. more precisely, the lenses $A$ and $B$ of transmissivity $a_{02}$ surrounded by an area $\Omega \backslash (A \cup B)$ of transmissivity $a_{01}$.



Figure 1

Then, the data $\bar{\mathbf{u}}$ are composed from these simulated measurements and – on every node where no measurement is given – from potential values $\tilde{\mathbf{u}} = (\mathbf{L}[\mathbf{a}_0])^{-1}\mathbf{d}$ resulting from the a priori guess $\mathbf{a}_0$. For technical reasons we set $\mu = \varrho^{-1}$.

*Results:* The Figures 2 to 5 show the transmissivity gained by

(1) direct inversion of the data $\bar{\mathbf{u}}$ (Fig. 2),

(2) inversion after data preparation, $\mu = 10^{-5}$ (Fig. 3),

(3) inversion after data preparation, $\mu = 5 \cdot 10^{-6}$ (Fig. 4),

(4) inversion after data preparation, $\mu = 10^{-6}$ (Fig.5).

*Discussion:* The lens $C$ in (1) - (4) is generally reproduced satisfactorily. This is expected since the data $\tilde{\mathbf{u}}$ are disturbed only in the area $B'$, i.e. disturbances in $B'$ do not essentially affect the area $C$. It confirms the above–mentioned local behavior of Vainikko's method. Moreover, Fig. 2 shows that the lens $B$ cannot be reconstructed by direct inversion of the data. Fortunately, this can be achieved by additional data preparation according to Section 3. The best reconstruction is obtained for $\mu = 10^{-6}$ (Fig. 5), where $B$ has nearly its correct shape, but its mean value is between $a_{01}$ and $a_{02}$, i.e. much less than its true value $a_{02}$. The latter fact is not surprising since the used (prepared) data are situated between $\tilde{\mathbf{u}}$ and $\bar{\mathbf{u}}$. The value $10^{-6}$ for $\mu$ seems to be optimal in this context but also $\mu = 5 \cdot 10^{-6}$ (Fig. 4) or $\mu = 5 \cdot 10^{-7}$ would be possible. But, for numerical reasons, $\mu \leq 10^{-7}$ is not suitable; in that case the condition number of the matrix $(\mu\mathbf{I} + \mathbf{L}[\mathbf{a}_0])^2$ appeared to be too bad for a calculation.
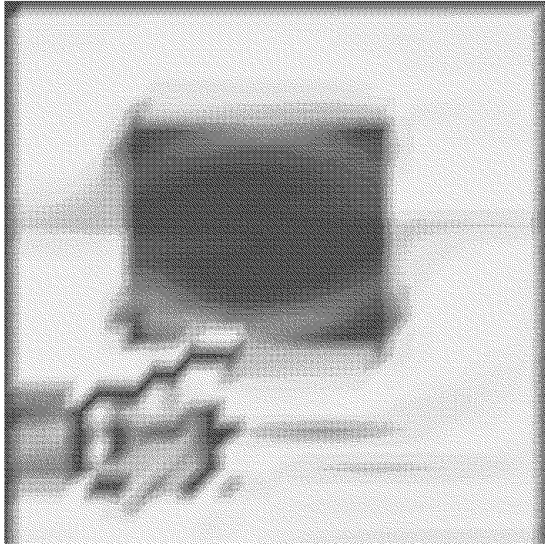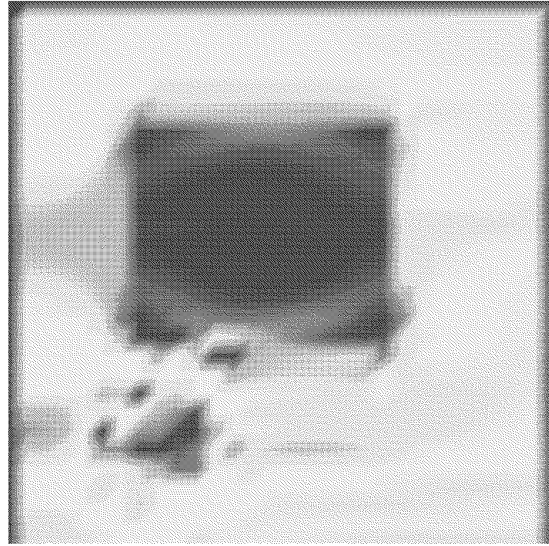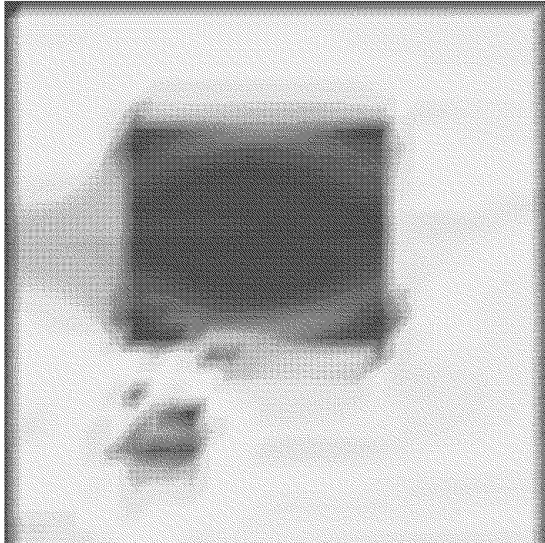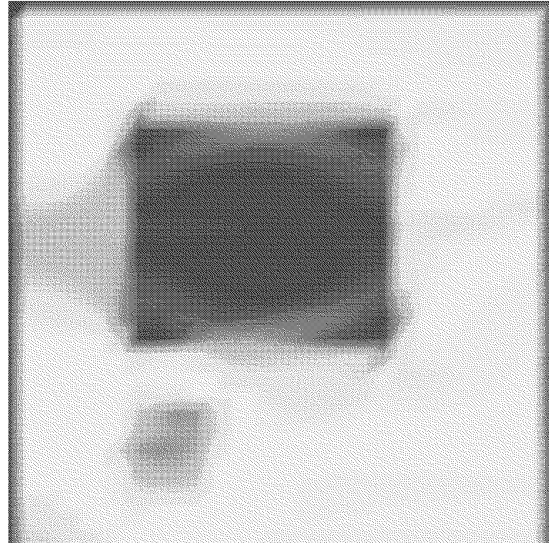
Figure 2



Figure 3



Figure 4



Fig. 5

In our example we have set $\mathbf{B} = \mathbf{I}$. We have made four calculations with five iterations each.

| Calculation | 1 | | 2 | | 3 | | 4 | |
|---|---|---|---|---|---|---|---|---|
| Iteration | $\mu$ | FIT | $\mu$ | FIT | $\mu$ | FIT | $\mu$ | FIT |
| 1 | $10^{-6}$ | 34948. | $10^{-6}$ | 34948. | $10^{-6}$ | 34948. | $10^{-6}$ | 34948. |
| 2 | $10^{-6}$ | 9467. | $2 \cdot 10^{-6}$ | 6377. | $2 \cdot 10^{-6}$ | 6377. | $10^{-5}$ | 1970. |
| 3 | $10^{-6}$ | 3944. | $3 \cdot 10^{-6}$ | 1736. | $4 \cdot 10^{-6}$ | 1433. | $10^{-4}$ | 33. |
| 4 | $10^{-6}$ | 2520. | $4 \cdot 10^{-6}$ | 785. | $8 \cdot 10^{-6}$ | 425. | $10^{-3}$ | 0.23 |
| 5 | $10^{-6}$ | 2038. | $5 \cdot 10^{-6}$ | 424. | $16 \cdot 10^{-6}$ | 126. | $10^{-2}$ | 0.002 |

with

$$\text{FIT} = \|\hat{\mathbf{u}} - \mathbf{b}\|_m^2.$$

*Description of the results*

**Calculation 1** ($\mu = 10^{-6}$ for every iteration). The area **B** being rather pale for one iteration now appears clearer and darker. This effect is clearly stronger than a visible increase of disturbances.

**Calculation 2 and 3** ($\mu = l \cdot 10^{-6}$, $l = 1, 2, 3, 4, 5$, and $\mu = l \cdot 10^{-6}$, $l = 1, 2, 2^2, 2^3, 2^4$), respectively. In principle we have the same as in calculation 1, but the contour of the area **B** is even clearer. Disturbances are tolerable but stronger.

**Calculation 4** ($\mu = 10^{-l}$, $l = 6, 5, 4, 3, 2$).
Here the rough structure is just discernible. But considerable disturbances have to be considered as destructive in a more complicated structure.

*Discussion.* 0 First of all it is clear that disturbances caused by uncertain data will increase by iteration. The reason is the ill–posedness of the inverse problem. Therefore, the following rule is useful

**Rule 1.** Iterate only a few times.

Another error source consists in the lack of measurements. If $\mu$ is large then $\hat{\mathbf{u}}$ is strongly influenced by the measured values. This can be disadvantageous in the inversion. The data $\hat{\mathbf{u}}$ is the more suitable the smaller $\mu$ is.

**Rule 2.** Choose $\mu$ small. But large enough that a difference to the a priori guess is visible. Then try to increase this trend by iteration.

As the test calculations show, the quantity FIT does not fit as a stopping rule. The reason is that there are a lot of **u** fitting the measurements, suitable and unsuitable ones.

**Rule 3.** Stop the iteration if the calculated $a$ does not change any more.

4.2. **A numerical example with real life data.** The following figures illustrate some numerical examples using real life data for an unconfined aquifer.

The flow region $\Omega$ of a 50 $km^2$ size was discretized by a finite element grid with 5365 nodes and 10479 triangles due to the details known from the geological point of view.

Inside $\Omega$ there are placed 50 observation points at which measured values of the groundwater level are available. Further values of the piezometric head are given at some ditches implemented by boundary conditions of third kind.

The Figure 6 describes the a priori guess for the transmissivity due to geological considerations. Figure 7 presents the reconstruction of the transmissivity from simulated measurements at each node. To simulate these measurements the a priori guess was used. In this case the parameter $\mu = 0$ and the regularization parameter $\alpha = 0$. Figure 8 shows the influence of the measurements mentioned above. The parameter $\mu$ can be understood as a weight between the influence of the measurements and of the a priori guess. We have chosen $\mu = 1$ and $\alpha = 0$. Finally, the Figure 9 demonstrates the smoothing properties of the Tikhonov regularization with choice of $\mu = 1$ and $\alpha = 10^{-2}$.
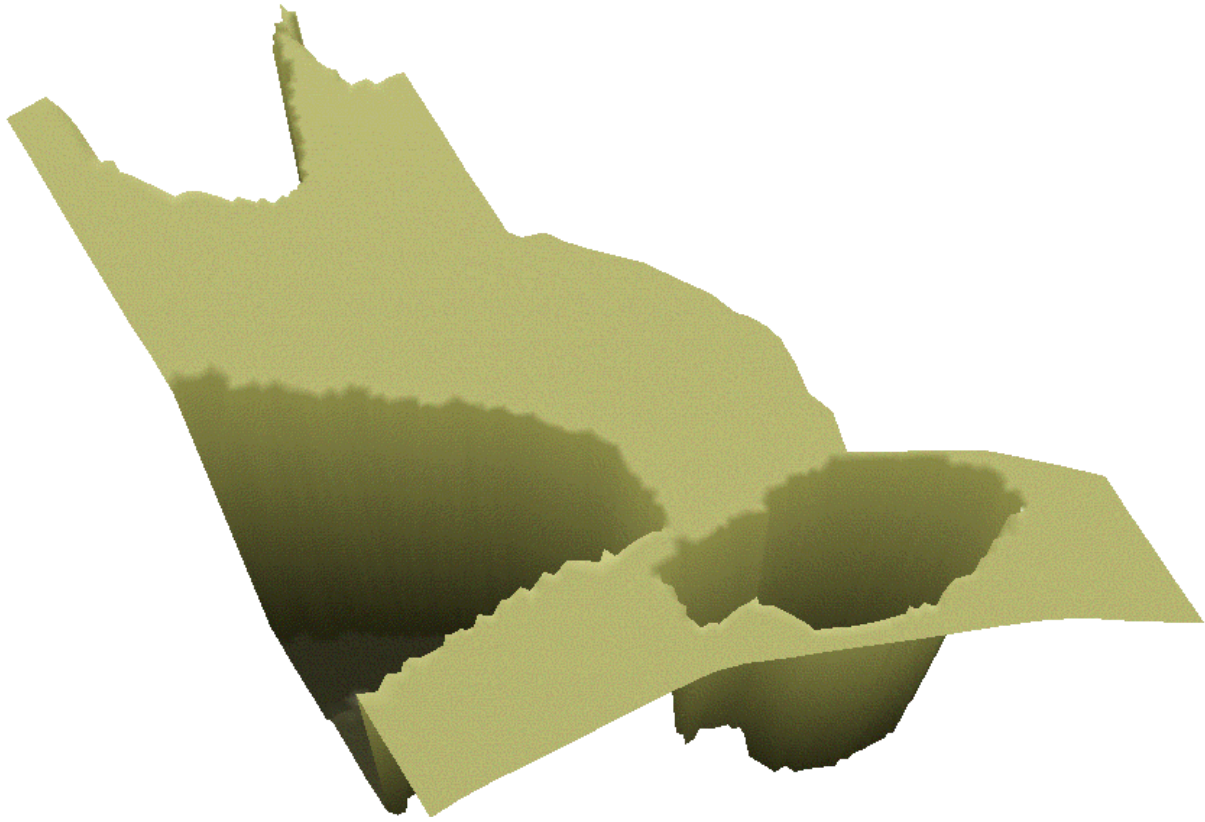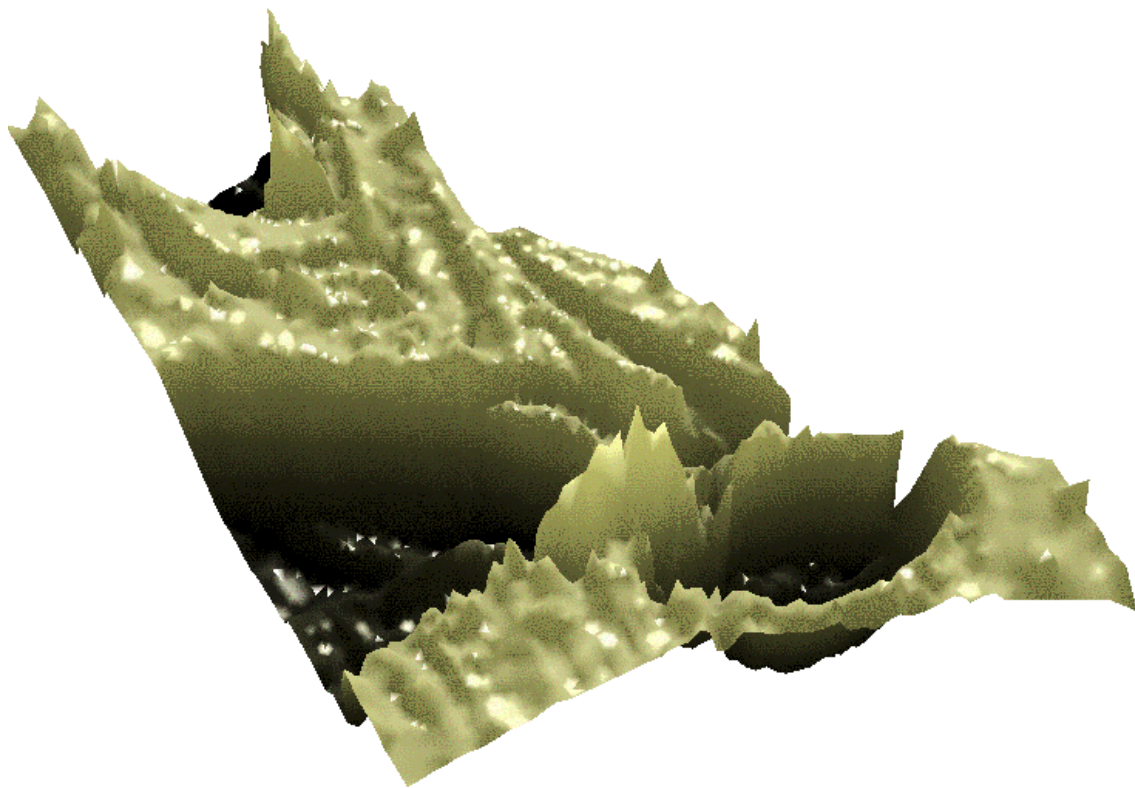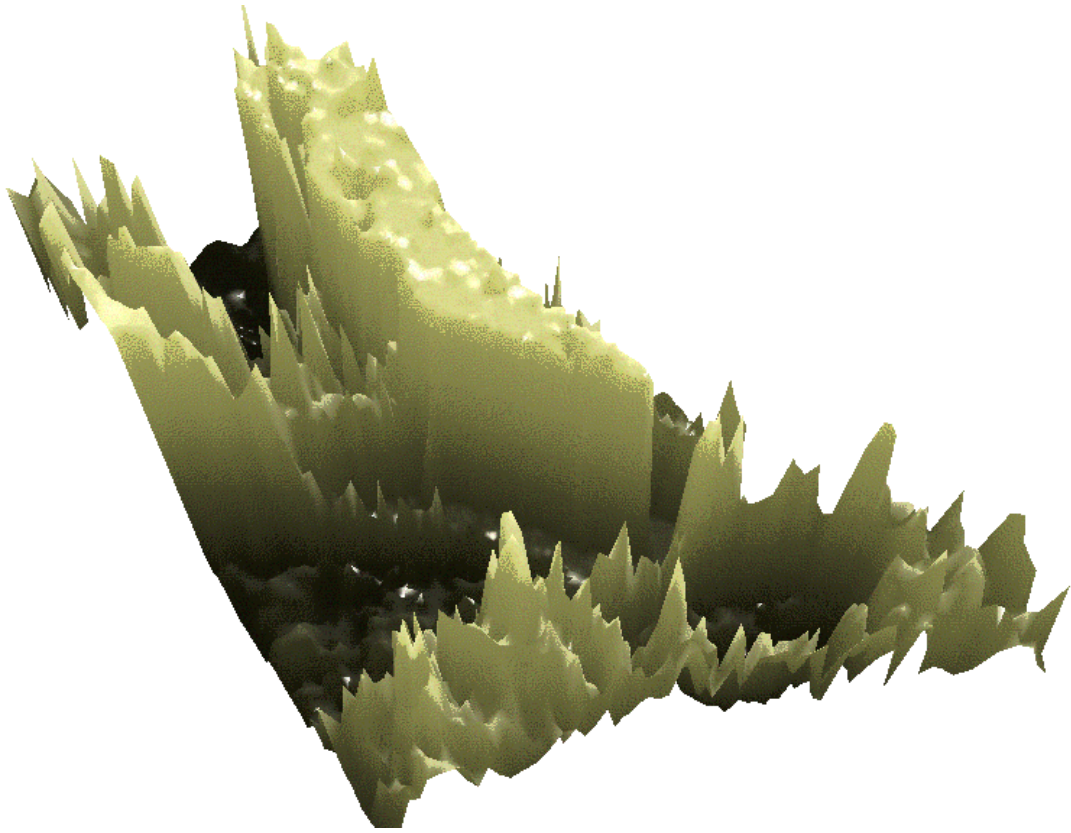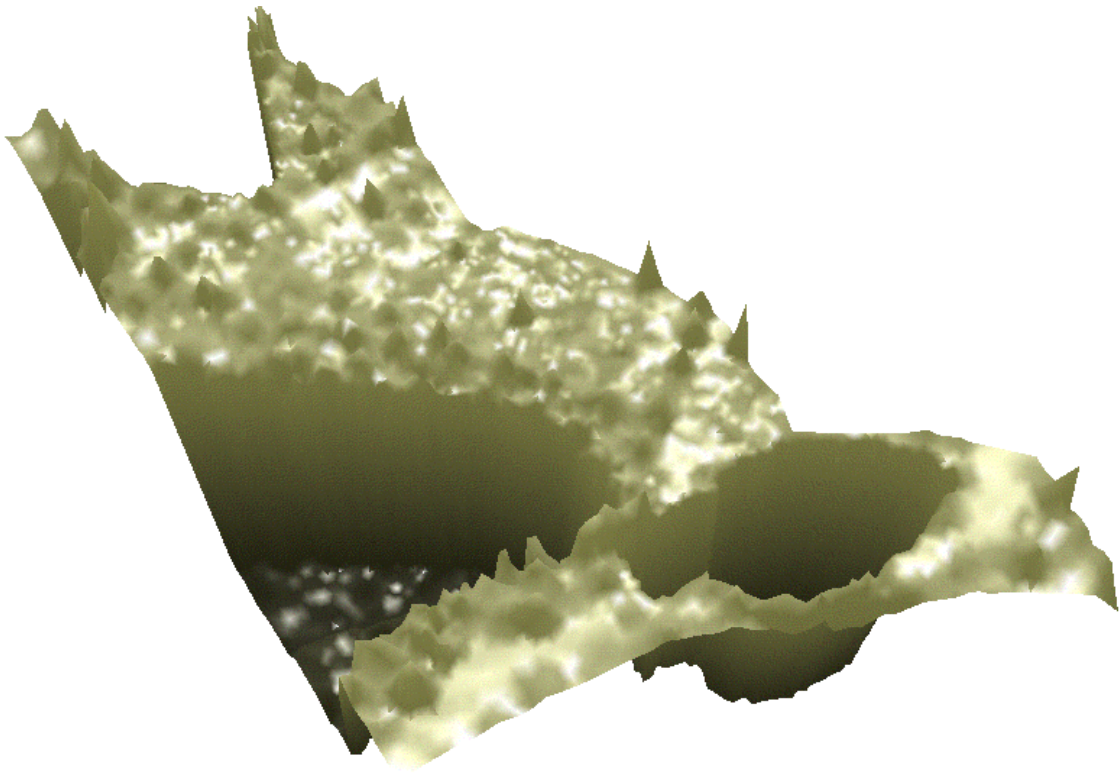
Figure 6



Figure 7

Figure 8



Figure 9

## References

[1] Acar, R.: Identification of the coefficient in elliptic equations, SIAM J. Control and Optimization **31** (1993), 1221–1244.

[2] Alessandrini, G.: An identification problem for an elliptic equation in two variables, Annali di Matematica Pura ed Applicata **145** (1986), 265–296.

[3] Bruckner, G., Handrock–Meyer, S. and Langmach, H.: On the identification of soil transmissivity from measurements of the groundwater level, WIAS–Preprint No. 250, Berlin 1996.

[4] Gottlieb, J. and Dietrich, P.: Identification of the permeability distribution in soil by hydraulic tomography, Inverse Problems **11** (1995), 353–360.

[5] Guidici, M.: Identifiability of distributed physical parameters in diffusive–like systems, Inverse problems **7** (1991), 231–245.

[6] Guidici, M., Morossi, G., Parravicini, G. and Ponzini, G.: A new method for the identification of distributed transmissivities, to appear in Water Resources Research.

[7] Hoffmann, K.–H. and Sprekels, J.: On the identification of elliptic problems by asymptotic regularization, Numer. Funct. Anal. and Optim. **7** (1984/85), 157–178.

[8] Hoffmann, K.–H. and Sprekels, J.: On the identification of parameters in general variational inequalities by asymptotic regularization, SIAM J. Math. Anal. **17** (1986), 1198–1217.

[9] Ito, K. and Kunisch, K.: A hybrid method combining the Output Least Squares and the Equation Error approach for the estimation of parameters in elliptic systems, SIAM J. Control Optim. **28** (1990), 113–136.

[10] Krein, S. G.: Linear equations in Banach spaces, Nauka, Moscow 1972.

[11] Ladyženskaya, O. A and Ural'ceva, N. N.: Linear and quasilinear elliptic equations, Academic Press, New York 1968.

[12] Lowe, B. and Kohn, R.V.: A variational method for numerically identifying a variable coefficient, in: Proceedings Int. Symp. on Variational Methods in the Geophysical Sciences, Norman, OK, USA 1985.

[13] Parker, R.L.: Geophysical inverse theory, Princeton University Press, Princeton 1994.

[14] Parravicini, G., Guidici, M., Morossi, G. and Ponzini, G.: Minimal apriori assignment in a direct method for determining phenomenological coefficients uniquely, to appear in Inverse Problems.

[15] Richter, G.R.: Numerical identification of a spatially varying diffusion coefficient, Math. Comp. **36** (1981), 375–386.

[16] Sprekels, J.: Identification of parameters in distributed systems: an overview, in: Methods of Operations Research **54**, 163–176, Verlag Anton Hain 1986.

[17] Sun, N.Z.: Inverse Problems in Groundwater Modeling, Kluwer Academic Publishers, Dordrecht 1994.

[18] Vainikko, G. and Kunisch, K.: Identifiability of the transmissivity coefficient in an elliptic boundary value problem, Zeitschrift für Analysis und ihre Anwendungen **12** (1993), 327–341.

[19] Vainikko, G.: Identification of filtration coefficient, A. Tikhonov (Ed.), Ill–Posed Problems in Natural Sciences (1992), 202–213.

[20] Vainikko, G.: On the discretization and regularization of ill–posed problems with noncompact operators, Numer. Funct. Anal. and Optim. **13** (1992), 381–396.